

Impact of fault prediction on checkpointing strategies

Guillaume Aupy¹, Yves Robert^{1,2}, Frédéric Vivien¹ and Dounia Zaidouni¹

1. Ecole Normale Supérieure de Lyon & INRIA, France

{Guillaume.Aupy|Yves.Robert|Frederic.Vivien|Dounia.Zaidouni}@ens-lyon.fr

2. University of Tennessee Knoxville, USA

July 31, 2012

Abstract

This paper deals with the impact of fault prediction techniques on checkpointing strategies. We extend the classical analysis of Young in the presence of a fault prediction system, which is characterized by its recall and its precision, and which provides either exact or window-based time predictions. We succeed in deriving the optimal value of the checkpointing period (thereby minimizing the waste of resource usage due to checkpoint overhead) in all scenarios. These results lay the foundations for future experimental validation of the model.

1 Introduction

In this paper, we assess the impact of fault prediction techniques on checkpointing strategies. We assume to have jobs executing on a platform subject to faults, and we let μ be the mean time between faults (MTBF) of the platform. In the absence of fault prediction, the standard approach is to take periodic checkpoints, each of length C , every period of duration T . In steady-state utilization of the platform, the value T_{opt} of T that minimizes the (expectation of the) waste of resource usage due to checkpointing, is easily computed as $T_{\text{opt}} = \sqrt{2C\mu}$. This is the well-known Young formula [4].

Now, when some fault prediction mechanism is available, can we compute a better checkpointing period to decrease the expected waste? and to what extent? Critical parameters that characterize a fault prediction system are its recall r , which is the fraction of faults that are indeed predicted, and its precision p , which is the fraction of predictions that are correct (i.e., correspond to actual faults). The major objective of this paper is to refine the expression of the expected waste as a function of these new parameters, and to derive optimal values for the checkpointing period. We deal with two problem instances, one where the predictor system provides exact dates for predicted events, and another where it only provides time windows during which events take place. We succeed in characterizing optimal values for both instances.

The results of this preliminary work lay the theoretical foundations for the study of the impact of prediction on checkpointing strategies, and are a prerequisite for conducting experimental simulations to fully validate the analysis for realistic application/platform scenarios.

2 Framework

2.1 Checkpointing strategy

We consider a *platform* subject to faults. Our work is agnostic of the granularity of the platform, which may consist either of a single processor, or of several processors that work concurrently and use coordinated checkpointing. The key parameter is μ , the mean time between faults (MTBF) of the platform. If the platform is made of K components whose individual MTBF is μ_{ind} , then $\mu = \frac{\mu_{\text{ind}}}{K}$.

Checkpoints are taken at regular intervals, or periods, of length T . We use C , D , and R for the duration of the checkpoint, downtime and recovery (respectively). We must enforce that $C \leq T$, and useful work is done only $T - C$ units of time during every period of length T , if no fault occurs. The

waste due to checkpointing in a fault-free execution is $\text{WASTE} = \frac{C}{T}$. In the following, the *waste* always denote the fraction of time that the platform is not doing useful work.

2.2 Fault predictor

A fault predictor is a mechanism that is able to predict that some faults will take place, either at a certain point in time, or within some time-interval window. The accuracy of the fault predictor is characterized by two quantities, the *recall* and the *precision*:

- The recall r is the fraction of faults that are predicted;
- The precision p is the fraction of fault predictions that are correct.

Traditionally, one defines three types of *events*: (i) *true positive* events are faults that the predictor has been able to predict (let True_P be their number); (ii) *false positive* events are fault predictions that did not materialize as actual faults (let False_P be their number); and (iii) *false negative* events are faults that were not predicted (let False_N be their number). With these definitions, we have

$$r = \frac{\text{True}_P}{\text{True}_P + \text{False}_N} \quad \text{and} \quad p = \frac{\text{True}_P}{\text{True}_P + \text{False}_P}$$

We point out that the precision p is a standard notion in the literature [1, 2, 3, 5, 6]. However, the name “precision” can be misleading. For instance, consider a predictor that provides time windows of length I for a platform whose MTBF is μ . The probability that a fault takes place inside the window is $\frac{I}{\mu}$, hence a precision $p = \frac{I}{\mu}$ brings no additional information. Of course a very high precision enables to identify those time intervals where faults are more likely to strike, but a very low precision is useful too (somewhat counter-intuitively!): it enables to identify those time intervals where faults should not be expected.

2.3 Fault rates

In addition to μ , the mean time between faults (MTBF) of the platform, let μ_P be the mean time between predicted events (both true positive and false positive), and let μ_{NP} be the mean time between unpredicted faults (false negative). Finally, we define the mean time between events as μ_e (including all three event types). The relationships between μ , μ_P , μ_{NP} , and μ_e are the following:

- $\frac{(1-r)}{\mu} = \frac{1}{\mu_{NP}}$ (here, $1 - r$ is the fraction of faults that are unpredicted);
- $\frac{r}{\mu} = \frac{p}{\mu_P}$ (here, r is the fraction of faults that are predicted, and p is the fraction of fault predictions that are correct);
- $\frac{1}{\mu_e} = \frac{1}{\mu_P} + \frac{1}{\mu_{NP}}$ (here, events are either predicted (true or false) or not).

3 Predictor with exact event dates

In this section, we present an analytical model to assess the impact of prediction on periodic checkpointing strategies. We consider the case where the predictor is able to provide exact prediction dates, and to generate such predictions at least C seconds in advance, so that a checkpoint can indeed be taken before the event (otherwise the prediction cannot be useful, because there is not enough time to take proactive actions). We consider the following algorithm:

1. While no fault prediction is available, checkpoints are taken periodically with period T .
2. When a fault is predicted, we decide whether to take the prediction into account or not. This decision is randomly taken: with probability q , we trust the predictor and take the prediction into account, and, with probability $1 - q$, we ignore the prediction. If we take the prediction into account, there are two cases. If we have enough time before the prediction date, we take a checkpoint as late as possible, i.e., so that it completes right at the time where the fault is predicted to happen. After the checkpoint, we then complete the execution of the period (see Figure 1). Otherwise, if we

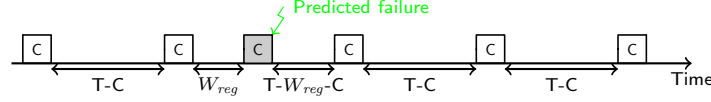


Figure 1: Whenever there is enough time, the algorithm takes a checkpoint just before the predicted failure.

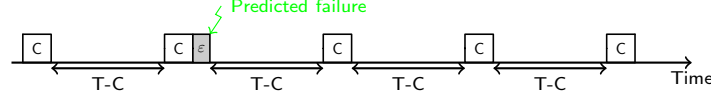


Figure 2: Whenever there is not enough time to take a checkpoint, the algorithm executes some extra work.

do not have enough time to take an extra checkpoint ($\varepsilon < C$), then we do some extra work during ε seconds (see Figure 2). We account for this work as idle time in the expression of the waste, to ease the analysis. Our expression of the waste is thus an upper bound.

The rationale for not always trusting the predictor is to avoid taking useless checkpoints too frequently. Intuitively, the precision p of the predictor must be above a given threshold for its usage to be worthwhile. In other words, if we decide to checkpoint just before a predicted event, there are two cases: either we will save time by avoiding a costly re-execution if the event does correspond to an actual fault, or we will lose time by unduly performing an extra checkpoint if the event does not correspond to an actual fault. We need a larger proportion of the former cases, i.e., a good precision, for the predictor to be really useful. The following analysis will determine the optimal value of q as a function of the parameters C , μ , r , and p .

3.1 Computing the waste

Our goal in this section is to compute a formula for the expected waste. Recall that the waste is the fraction of time that the processors do not perform useful computations, either because they are checkpointing, or because they recover from a fault. There are four different sources of waste (see Figure 3):

1. **Checkpoints:** During a fault-free execution, the fraction of resources used in checkpointing is:

$$\frac{C}{T} \quad (1a)$$

2. **Unpredicted faults:** This overhead occurs each time a unpredicted fault strikes, that is, on average, once every μ_{NP} seconds. The time wasted because of the unpredicted fault is then the time elapsed between the last checkpoint and the fault, plus the downtime and the time needed for the recovery. The expectation of the time elapsed between the last checkpoint and the fault is equal to half the period of checkpoints, because the time where the fault hits the system is independent of the checkpointing algorithm. Finally, the waste due to unpredicted faults is:

$$\frac{1}{\mu_{NP}} \left[\frac{T}{2} + D + R \right] \quad (1b)$$

3. **Predictions taken into account:** Now we have to compute the execution overhead due to a prediction which we trust (hence we checkpoint just before its date). This overhead occurs each time a prediction is made by the predictor, that is, on average, once every μ_P seconds, and that we decide to trust it, with probability q . If the predicted event is an actual fault, we waste $C + D + R$ seconds (we waste $D + R$ seconds because the predicted event corresponds to an actual fault and if we have enough time before the prediction date, we waste C seconds because we take an extra checkpoint as late as possible before the prediction date (see Figure 1). Note that if we do not have

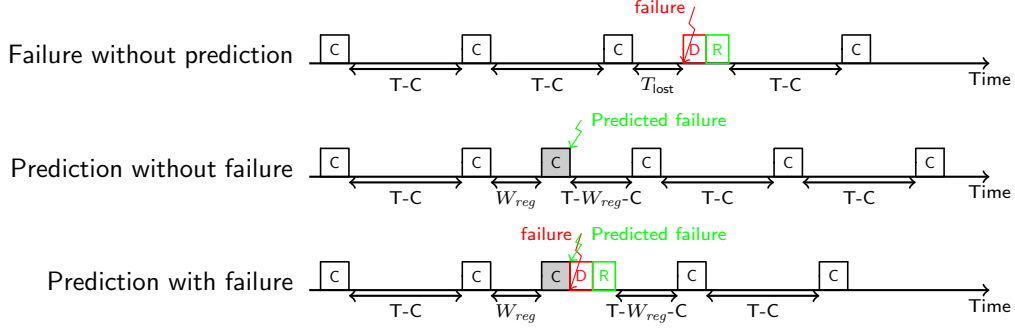


Figure 3: Actions taken when the predictor provides exact event dates.

enough time to take an extra checkpoint (see Figure 2), we overestimate the waste as C seconds. Otherwise, if the predicted event is not an actual fault, we waste C seconds. An actual fault occurs with probability p , and a false prediction is made with probability $(1 - p)$. Averaging with these probabilities, we waste an expected amount of $[p(C + D + R) + (1 - p)C]$ seconds. Finally, the corresponding overhead is:

$$\frac{1}{\mu_P} q [p(C + D + R) + (1 - p)C] \quad (1c)$$

4. **Ignored predictions:** The final source of waste is for predicted events that we do not trust. This overhead occurs each time a prediction is made by the predictor, that is, on average, once every μ_P seconds, and that we decide to trust it, with probability $1 - q$. If the predicted event corresponds to an actual fault, we waste $(\frac{T}{2} + D + R)$ seconds (as for a unpredicted fault). Otherwise there is no fault and we took no extra checkpoint, and thus we lose nothing. An actual fault occurs with a probability p . The corresponding overhead is:

$$\frac{1}{\mu_P} (1 - q) \left[p \left(\frac{T}{2} + D + R \right) + (1 - p)0 \right] \quad (1d)$$

Summing up the overhead over the four different sources, we obtain the following equation for the waste:

$$\begin{aligned} \text{WASTE} &= \frac{C}{T} \\ &+ \frac{1}{\mu_{NP}} \left[\frac{T}{2} + D + R \right] \\ &+ \frac{1}{\mu_P} q [p(C + D + R) + (1 - p)C] \\ &+ \frac{1}{\mu_P} (1 - q) \left[p \left(\frac{T}{2} + D + R \right) + (1 - p)0 \right] \end{aligned}$$

After simplification, we have:

$$\text{WASTE} = \frac{C}{T} + \frac{1}{\mu} \left[(1 - rq) \frac{T}{2} + D + R + \frac{qr}{p} C \right] \quad (2)$$

3.2 Validity of the analysis

We point out that Equation (2) is accurate only when two events (an event being a prediction (true or false) or a unpredicted fault) do not take place within the same period. To ensure that this condition is met with a high probability, we bound the length of the period: we enforce the condition $T < \alpha \mu_e$, where α is some tuning parameter.

In fact, the number of events during a period of length T can be modeled as a Poisson process of parameter $\frac{T}{\mu_e}$; the probability of having $k \geq 0$ faults is $\frac{1}{k!}(\frac{T}{\mu_e})^k e^{-\frac{T}{\mu_e}}$. Hence the probability of having two or more faults is $\pi = P(X \geq 2) = 1 - (P(X = 0) + P(X = 1)) = 1 - (1 + \frac{T}{\mu_e})e^{-\frac{T}{\mu_e}}$, where X is the number of faults. Enforcing the constraint $T \leq \alpha\mu_e$ leads to $\pi \leq 1 - (1 + \alpha)e^{-\alpha}$. If we assume $\alpha = 0.1$ then $\pi \leq 0.005$, hence a valid approximation when bounding the period range accordingly. Indeed, with such a conservative value for α , we have overlapping faults every 200 periods in average, so that the model is accurate for 98% of the checkpointing segments, hence quite reliable.

In addition to the previous constraint, recall that we must always enforce the condition $C \leq T$, by construction of the periodic checkpointing policy. Finally, note that the optimal waste may never exceed 1, since it represents the fraction of time that is “wasted”. When the waste equals 1, the application no longer makes progress.

3.3 Waste minimization

3.3.1 Computation of the extremum period T_{extr}

We have the expression of the waste from Equation (2):

$$\text{WASTE}(T) = \frac{C}{T} + \frac{1}{\mu} \left[(1 - rq) \frac{T}{2} + D + R + \frac{qr}{p} C \right]$$

Differentiating twice with respect to T , we obtain

$$\text{WASTE}'(T) = \frac{-C}{T^2} + \frac{1}{\mu} \left[(1 - rq) \frac{1}{2} \right]$$

$$\text{WASTE}''(T) = \frac{2C}{T^3} > 0$$

We obtain that $\text{WASTE}''(T)$ is strictly positive, hence $\text{WASTE}(T)$ is a convex function of T and admits a unique minimum on its domain $[C, \alpha\mu_e]$. We also compute T_{extr} , the extremum value of T that is the unique zero of the function $\text{WASTE}'(T)$:

$$T_{\text{extr}} = \sqrt{\frac{2\mu C}{1 - rq}}$$

Note that this Equation makes sense even when $1 - rq = 0$: indeed this would mean that both $r = 1$ and $q = 1$: the predictor predicts every fault, and we take proactive action for each one of them. Then there should never be any periodic checkpointing!

Finally, note that T_{extr} may well not belong to the admissible domain $[C, \alpha\mu]$. The optimal waste $\text{WASTE}_{\text{opt}}$ is determined via the following case analysis.

3.3.2 Computation of $\text{Waste}_{\text{opt}}$

We rewrite the waste as an affine function of q :

$$\text{WASTE}(q) = \frac{rq}{\mu} \left(\frac{C}{p} - \frac{T}{2} \right) + \left(\frac{C}{T} + \frac{T}{2\mu} + \frac{D + R}{\mu} \right)$$

For any value of T , we deduce that $\text{WASTE}(q)$ is minimized either for $q = 0$ or for $q = 1$. This (somewhat unexpected) conclusion is that the predictor should sometimes be always trusted, and sometimes never, but no in-between value for q will do a better job. Thus we need to minimize the two functions $\text{WASTE}_{\{q=0\}}$ and $\text{WASTE}_{\{q=1\}}$ over the domain of admissible values for T , and to retain the best result.

We have:

$$\text{WASTE}_{\{q=0\}}(T) = \frac{C}{T} + \frac{1}{\mu} \left[\frac{T}{2} + D + R \right]$$

The function $\text{WASTE}_{\{q=0\}}(T)$ is a convex function and reaches its minimum for T_{opt1} in the interval $[C, \alpha\mu]$:

- If $(C < T_{\text{extr}} < \alpha\mu_e)$: $T_{\text{opt1}} = T_{\text{extr}} = \sqrt{2\mu C}$
- If $(T_{\text{extr}} < C)$: $T_{\text{opt1}} = C$
- If $(T_{\text{extr}} \geq \alpha\mu_e)$: $T_{\text{opt1}} = \alpha\mu_e$

Thus, $\text{WASTE}_{\{q=1\}}$ is minimized for:

$$T_{\text{opt1}} = \min \left\{ \alpha\mu_e, \max\{\sqrt{2\mu C}, C\} \right\}$$

Similarly, we have:

$$\text{WASTE}_{\{q=1\}}(T) = \frac{C}{T} + \frac{1}{\mu} \left[(1-r)\frac{T}{2} + D + R + \frac{r}{p}C \right]$$

The function $\text{WASTE}_{\{q=1\}}(T)$ is a convex function and reaches its minimum for T_{opt2} in the interval $[C, \alpha\mu_e]$.

- If $(C < T_{\text{extr}} < \alpha\mu_e)$: $T_{\text{opt2}} = T_{\text{extr}} = \sqrt{\frac{2\mu C}{1-r}}$
- If $(T_{\text{extr}} < C)$: $T_{\text{opt2}} = C$
- If $(T_{\text{extr}} \geq \alpha\mu_e)$: $T_{\text{opt2}} = \alpha\mu_e$

Thus, $\text{WASTE}_{\{q=1\}}$ is minimized for:

$$T_{\text{opt2}} = \min \left\{ \alpha\mu_e, \max\left\{ \sqrt{\frac{2\mu C}{1-r}}, C \right\} \right\}$$

Finally, the optimal waste is:

$$\text{WASTE}_{\text{opt}} = \min \left\{ \text{WASTE}_{\{q=0\}}(T_{\text{opt1}}), \text{WASTE}_{\{q=1\}}(T_{\text{opt2}}) \right\}$$

3.4 Prediction and preventive migration

In this section, we make a short digression and briefly present an analytical model to assess the impact of prediction and preventive migration on periodic checkpointing strategies. As before, we consider a predictor that is able to predict exactly when faults happen, and to generate these predictions at least C seconds before the event dates.

The idea of migration consists in moving a task for execution on another node, when a fault is predicted to happen on the current node in the near future. Note that the faulty node can later be replaced, in case of a hardware fault, or software rejuvenation can be used in case of a software fault. We consider the following algorithm, which is very similar to that used in Section 3.1:

1. When no fault prediction is available, checkpoints are taken periodically with period T .
2. When a fault is predicted, we decide whether to execute the migration or not. The decision is a random one: with probability q we trust the predictor and do the migration and, with probability $1-q$, we ignore the prediction. If we take the prediction into account, we execute the migration as late as possible, so that it completes right at the time when the fault is predicted to happen.

As before, we have four different sources of waste. Summing the overhead of the execution of these different sources, we obtain the following equation for the waste (where M is the duration of a migration):

$$\begin{aligned} \text{WASTE} = & \frac{C}{T} \\ & + \frac{1}{\mu_{NP}} \left[\frac{T}{2} + D + R \right] \\ & + \frac{1}{\mu_P} q [p(M) + (1-p)M] \\ & + \frac{1}{\mu_P} (1-q) \left[p\left(\frac{T}{2} + D + R\right) + (1-p)0 \right] \end{aligned}$$

After simplification, we get:

$$\text{WASTE} = \frac{C}{T} + \frac{1}{\mu} \left[(1 - rq) \left(\frac{T}{2} + D + R \right) + \frac{qr}{p} M \right] \quad (4)$$

Equation (4) is very similar to Equation (2), and the minimization of the waste proceeds exactly as in Section 3.3. In a nutshell, $\text{WASTE}(T)$ is again a convex function and admits a unique minimum over its domain $[C, \alpha\mu_e]$, the unique zero of the derivative has the same value $T_{\text{extr}} = \sqrt{\frac{2\mu C}{1-rq}}$, and for any value of T , the waste is minimized for either $q = 0$ or $q = 1$. We conduct the very same case analysis as in Section 3.3.

4 Predictor with a prediction window

In the previous section, we have supposed that the predictor was able to predict exactly when faults will strike. Here, we suppose (maybe more realistically) that the predictor gives a *prediction window*, that is an interval of time of length I during which the predicted fault is likely to happen. As before in Section 3, we suppose that we have enough time to checkpoint before the beginning of the prediction window. Also, as in Section 3, when a prediction is made, we enforce that the scheduling algorithm has the choice either to take or not to take this prediction into account, with probability q .

We start with a description of the strategies that can be used, depending upon the (relative) length I of the prediction window. Let us define two *modes* for the scheduling algorithm:

Regular: This is the mode used when no fault prediction is available, or when a prediction is available but we decide to ignore it (with probability $1 - q$). In regular mode, we use periodic checkpointing with period T_{NP} . Intuitively, T_{NP} corresponds to the checkpointing period T of Section 3.

Proactive: This is the mode used when a fault prediction is available and we decide to trust it – decision taken with probability q –. Consider such a trusted prediction made with the prediction window $[t_0, t_0 + I]$. There are several strategies that can be envisioned:

1. *Instantaneous*– The first strategy is to ignore the time-window and to execute the same algorithm as if the predictor had given an exact fault date at time t_0 . Just as described in Section 3, the algorithm interrupts the current period (of scheduled length T_{NP}), checkpoints during the interval $[t_0 - C, t_0]$, and then returns to regular mode: at time t_0 , it resumes the work due to complete the interrupted period.
2. *No checkpoint during prediction window*– The second strategy is intended for a short prediction window: instead of ignoring it, we acknowledge it, but make the decision not to checkpoint during it. As in the first strategy, the algorithm interrupts the current period (of scheduled length T_{NP}), and checkpoints during the interval $[t_0 - C, t_0]$. But here, we return to regular mode only at time $t_0 + I$, where we resume the work due to complete the interrupted period of the regular mode. During the whole length of the time-window, we execute work without checkpointing, at the risk of losing work if a fault indeed strikes. But for a small value of I , it may not be worthwhile to checkpoint during the prediction window (if at all possible, note that there is no choice if $I < C$).
3. *With checkpoints during prediction window*– The third strategy is intended for a longer prediction window: as before, the algorithm interrupts the current period (of scheduled length T_{NP}), and checkpoints during the interval $[t_0 - C, t_0]$, but then decides to take several checkpoints during the prediction window. The period T_{P} of these checkpoints in proactive mode will presumably be shorter than T_{NP} , to take into account the higher fault probability. To simplify the presentation, we will use an integer number of periods of length T_{P} within the prediction window. In the following, we analytically compute the optimal number of such periods. But we take at least one period here, hence one checkpoint, which implies that $C \leq I$. We return to regular mode either right after the fault stroke within the time window $[t_0, t_0 + I]$, or at time $t_0 + I$ whenever no actual fault happens within the prediction window. Then, we resume the work due to complete the interrupted period of the regular mode.

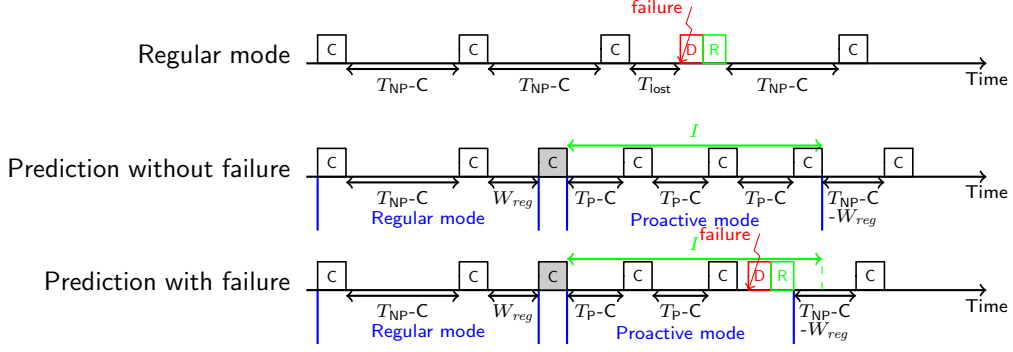


Figure 4: Outline of the behavior of Algorithm 1 (third strategy) (checkpoints taken during the prediction window in proactive mode).

The third strategy is the most complex to describe, and the complete behavior of the scheduling algorithm is shown in Algorithm 1. Note that for all strategies, exactly as in Section 3, we insert some additional work for the particular case where there is not enough time to take a checkpoint before entering proactive mode (because a checkpoint for the regular mode is currently on-going, see Figure 2). We account for this work as idle time in the expression of the waste, to ease the analysis. Our expression of the waste is thus an upper bound.

Algorithm 1: Proactive algorithm.

```

1 if fault happens then
2   | After downtime, execute recovery;
3   | Enter regular mode;
4 if in proactive mode for a time greater than or equal to I then
5   | Switch to regular mode
6 if Prediction made with interval  $[t, t + I]$  and prediction taken into account then
7   | Let  $t_C$  be the date of the last checkpoint under regular mode to start no later than  $t - C$ ;
8   | if  $t_C + C < t - C$  then (time for an extra checkpoint)
9   |   | Take a checkpoint starting at time  $t - C$ 
10  | else (no time for the extra checkpoint)
11  |   | Work in the time interval  $[t_C + C, t]$ 
12  |    $W_{reg} \leftarrow \max(0, t - C - (t_C + C))$ ;
13  |   Switch to proactive mode at time  $t$ ;
14 while in regular mode and no predictions are made and no faults happen do
15  |   Work for a time  $T_{NP} - W_{reg} - C$  and then checkpoint;
16  |    $W_{reg} \leftarrow 0$ ;
17 while in proactive mode and no faults happen do
18  |   Work for a time  $T_P - C$  and then checkpoint;
```

First we compute the waste incurred by the three algorithms, starting with the most complex strategy (Section 4.1), and then simplifying the formula and establishing the result for the other two strategies (Section 4.2). Then we discuss the validity of the model in Section 4.3. Finally, we solve the optimization problem and derive optimal values for the parameter q , and for the two periods T_P and T_{NP} (Section 4.4).

4.1 Computing the waste with checkpoints during prediction window

In this section we focus on computing the waste of the most complex strategy, that with checkpoints during prediction window (Algorithm 1). As in Section 3, we assume that there is a single event of any type (either a prediction (true or false), or an unpredicted failure). As already mentioned, we discuss this hypothesis in Section 4.3.

We first compute which fraction of the time the algorithm spends in either mode:

- the fraction of time spent in the *regular* mode (checkpointing with period T_{NP});
- the fraction of time spent in the *proactive* mode (checkpointing with period T_P).

Let I' be the average time spent in the *proactive* mode. When a prediction is made, we may choose to ignore it, which happens with probability $1 - q$. In this case, the algorithm stays in regular mode and does not spend any time in the proactive mode. With probability q , we may decide to take the prediction into account. In this case, if the prediction is a false positive event (no actual fault strikes), which happens with probability $1 - p$, then the algorithm spends I units of time in the proactive mode. Otherwise, if the prediction is a true positive event (an actual fault hits the system), which happens with probability p , then the algorithm spends an average of $\mathbb{E}_I^{(f)}$ in the proactive mode. Here $\mathbb{E}_I^{(f)}$ is the expectation of the time elapsed between the beginning of the prediction window and the time when a fault happens, knowing that a fault happens in the prediction window. Note that if faults are uniformly distributed across the prediction window, then $\mathbb{E}_I^{(f)} = \frac{I}{2}$. Altogether, we obtain:

$$I' = (1 - q) \cdot 0 + q \left((1 - p) \cdot I + p \cdot \mathbb{E}_I^{(f)} \right) = q \left((1 - p)I + p\mathbb{E}_I^{(f)} \right).$$

Finally, each time there is a prediction, that is, on the average, every μ_P seconds, the algorithm spends a time I' in the proactive mode. Therefore, Algorithm 1 spends a fraction of time $\frac{I'}{\mu_P}$ in the proactive mode, and a fraction of time $1 - \frac{I'}{\mu_P}$ in the regular mode. We now identify the four different sources of waste, and we analyze their respective costs.

1. **Waste due to periodic checkpointing.** There are two cases, depending upon the mode of Algorithm 1.

- (a) **Regular mode.** In this mode, we take periodic checkpoints. We take a checkpoint of size C each time the algorithm has processed work for a time $T_{NP} - C$ in the regular mode. This remains true if, after spending some time in the regular mode, the algorithm switches to the proactive mode, and later switches back to the regular mode. This behavior is enforced by recording the amount of work performed under the regular mode (variable W_{reg} , at line 12 of Algorithm 1), and by taking this value into account at line 15.

Taking into account the fraction of time that Algorithm 1 spends in the regular mode, this source of waste has a total cost of:

$$\left(1 - \frac{I'}{\mu_P} \right) \frac{C}{T_{NP}}. \quad (5a)$$

- (b) **Proactive mode.** In this mode, we take a checkpoint of size C each time the algorithm has processed work for a time $T_P - C$.

If no fault happens while the algorithm is in the proactive mode, then the algorithm stays exactly a time I in this mode (thanks to the condition at line 4). The waste due to the periodic checkpointing is exactly $\frac{C}{T_P}$ (because T_P divides I).

If a fault happens while the algorithm is in proactive mode, then, the expectation of the waste due to the periodic checkpointing is upper-bounded by the same quantity $\frac{C}{T_P}$. This is an over-approximation of the waste in that case, because the fault may strike before full completion of the last period.

Overall, taking into account the fraction of time Algorithm 1 is in the proactive mode, the cost of this source of waste is:

$$\frac{I'}{\mu_P} \frac{C}{T_P}. \quad (5b)$$

2. **Waste incurred when switching to the proactive mode.** Each time we take into account a prediction (which happens with probability q on average every μ_P units of time), we start by doing one preliminary checkpoint if we have the time to do so (line 9). If we do not have the time to take an additional checkpoint, the algorithm do not do any processing for a duration of at most C (line 11). In both cases, the wasted time is at most C and this happens once every $\frac{\mu_P}{q}$. Hence, switching from the regular mode to the proactive one induces a waste of at most

$$\frac{q}{\mu_P} C \quad (5c)$$

3. **Waste due to predicted faults.** Predicted faults happen with frequency $\frac{p}{\mu_P}$. As we may choose to ignore a prediction, there are still two cases depending on the mode of the algorithm at the time the fault hits the system.

(a) **Regular mode.** If the algorithm is in regular mode when a predicted fault hits, this means that we have chosen to ignore the prediction, a decision taken with probability $(1 - q)$.

The time wasted because of the predicted fault is then the time elapsed between the last checkpoint and the fault, plus the downtime and the time needed for the recovery. The expectation of the time elapsed between the last checkpoint and the fault is equal to half the period of checkpoints, because the time where the fault hits the system is independent of the checkpointing algorithm. Therefore, the waste due to predicted faults hitting the system in regular mode is:

$$\frac{p(1-q)}{\mu_P} \left(\frac{T_{NP}}{2} + D + R \right) \quad (5d)$$

(b) **Proactive mode.** If the algorithm is in proactive mode when a fault hits, then we have chosen to take the prediction into account, a decision that is taken with probability q .

The time wasted because of the predicted fault is then, in addition to the downtime and the time needed for the recovery, the time elapsed between the last checkpoint and the fault or, if no checkpoint had already been taken in the proactive mode, the time elapsed between the start of the proactive mode and the fault.

Here, we can no longer assume that the time the fault hits the system is independent of the checkpointing date. This is because the proactive mode starts exactly at the beginning of the prediction window. Let T_{lost} denote the computation time elapsed between the latest of the beginning of the proactive mode and the last checkpoint, and the fault date. Then the expectation of T_{lost} depends on the distribution of the fault date in the prediction window. However, we know that whatever the distribution, $T_{lost} \leq T_P$. Therefore we over approximate the waste in that case by:

$$\frac{qp}{\mu_P} (T_P + D + R) \quad (5e)$$

4. **Waste due to unpredicted faults.** There are again two cases, depending upon the mode of the algorithm at the time the fault hits the system.

(a) **Regular mode.** In this mode the work done is periodically checkpointed with period T_{NP} .

The time wasted because of an unpredicted fault is then the time elapsed between the last checkpoint and the fault, plus the downtime and the time needed for the recovery. As before, the expectation of this value is $T_{lost} = \frac{T_{NP}}{2}$.

An unexpected fault hits the system once every μ_{NP} seconds on the average. Taking into account the fraction of the time the algorithm is in regular mode, the waste due to unpredicted faults hitting the system in regular mode is:

$$\left(1 - \frac{I'}{\mu_P} \right) \frac{1}{\mu_{NP}} \left(\frac{T_{NP}}{2} + D + R \right) \quad (5f)$$

(b) **Proactive mode.** Because of the assumption that a single event takes place within a time-interval, we do not consider the very unlikely case where a unpredicted fault strikes during a prediction window. This amounts to assume that $\frac{I'}{\mu_P} \frac{1}{\mu_{NP}} (T_P + D + R)$ is negligible.

Finally, we gather the expressions of the different types of waste of Equations (5a) through (5f) to obtain the formula of the overall waste:

$$\begin{aligned} \text{WASTE}_{\text{withCkpt}} &= \left(1 - \frac{I'}{\mu_P} \right) \frac{C}{T_{NP}} + \frac{I'}{\mu_P} \frac{C}{T_P} + \frac{q}{\mu_P} C + \frac{p(1-q)}{\mu_P} \left(\frac{T_{NP}}{2} + D + R \right) + \frac{pq}{\mu_P} (T_P + D + R) \\ &\quad + \left(1 - \frac{I'}{\mu_P} \right) \frac{1}{\mu_{NP}} \left(\frac{T_{NP}}{2} + D + R \right) \\ \text{WASTE}_{\text{withCkpt}} &= \left(\left(1 - \frac{I'}{\mu_P} \right) \frac{1}{T_{NP}} + \frac{I'}{\mu_P} \frac{1}{T_P} + \frac{q}{\mu_P} \right) C + \frac{p(1-q)}{\mu_P} \frac{T_{NP}}{2} + \frac{pq}{\mu_P} T_P + \left(1 - \frac{I'}{\mu_P} \right) \frac{1}{\mu_{NP}} \frac{T_{NP}}{2} \\ &\quad + \left(\frac{p}{\mu_P} + \left(1 - \frac{I'}{\mu_P} \right) \frac{1}{\mu_{NP}} \right) (D + R) \end{aligned} \quad (6)$$

4.2 Computing the waste of the other strategies

The waste of the first strategy (*Instantaneous*) is very close to the one given in Equation (2). The difference lies in T_{lost} , the expectation of the work lost when a fault is predicted and the prediction is taken into account. When a prediction is taken into account and the predicted event is an actual fault, the waste in Equation (2) was $\frac{qp}{\mu_P}(C + D + R)$ (see Equation (1c)). Because the prediction was exact, T_{lost} was equal to 0. However in our new Equation, the waste for this part is now $\frac{qp}{\mu_P}(C + T_{\text{lost}} + D + R)$. On average, the fault occurs after a time $\mathbb{E}_I^{(f)}$. However, because we do not know the relation between $\mathbb{E}_I^{(f)}$ and T_{NP} , then T_{lost} has expectation $\frac{T_{\text{NP}}}{2}$ if $\frac{T_{\text{NP}}}{2} \leq \mathbb{E}_I^{(f)}$. The new waste is then:

$$\text{WASTE}_{\text{instant}} = \frac{C}{T_{\text{NP}}} + \frac{1}{\mu} \left[(1 - rq) \frac{T_{\text{NP}}}{2} + D + R + \frac{qr}{p} C + qr \min \left(\mathbb{E}_I^{(f)}, \frac{T_{\text{NP}}}{2} \right) \right] \quad (7)$$

As for the second strategy (*No checkpoint during prediction window*), we do no longer incur the waste of Equation (5b) as we no longer checkpoint in proactive mode. Furthermore, the value of T_{lost} in Equation (5e) becomes $\mathbb{E}_I^{(f)}$ instead of T_P . Consequently, the total waste when there is no checkpoint during the proactive mode is:

$$\begin{aligned} \text{WASTE}_{\text{noCkpt}} &= \left(1 - \frac{I'}{\mu_P} \right) \frac{C}{T_{\text{NP}}} + \frac{q}{\mu_P} C + \frac{p(1-q)}{\mu_P} \left(\frac{T_{\text{NP}}}{2} + D + R \right) + \frac{pq}{\mu_P} \left(\mathbb{E}_I^{(f)} + D + R \right) \\ &\quad + \left(1 - \frac{I'}{\mu_P} \right) \frac{1}{\mu_{\text{NP}}} \left(\frac{T_{\text{NP}}}{2} + D + R \right) \end{aligned}$$

which we rewrite as

$$\begin{aligned} \text{WASTE}_{\text{noCkpt}} &= \left(\left(1 - \frac{I'}{\mu_P} \right) \frac{1}{T_{\text{NP}}} + \frac{q}{\mu_P} \right) C + \frac{p(1-q)}{\mu_P} \frac{T_{\text{NP}}}{2} + \frac{pq}{\mu_P} \mathbb{E}_I^{(f)} + \left(1 - \frac{I'}{\mu_P} \right) \frac{1}{\mu_{\text{NP}}} \frac{T_{\text{NP}}}{2} \\ &\quad + \left(\frac{p}{\mu_P} + \left(1 - \frac{I'}{\mu_P} \right) \frac{1}{\mu_{\text{NP}}} \right) (D + R) \end{aligned} \quad (8)$$

Note that when $I = 0$, the first and second strategies collapse. Indeed, we have $\mathbb{E}_I^{(f)} = 0$ if $I = 0$, and we check that Equations (7) and (8) are identical in that case.

4.3 Validity of the results

In this subsection, we discuss the validity of the model. The analysis is similar to that of Section 3.2, except that we deal with different length intervals here. As before, we assume that there is a single event of any type within each interval under study. The condition $T \leq \alpha\mu_e$ then becomes

$$T_{\text{NP}} + I \leq \alpha\mu_e \quad (9)$$

as $T_{\text{NP}} + I$ is the longer time interval considered in the analysis of Algorithm 1.

4.4 Waste minimization

In this section we aim at minimizing the waste of the three strategies, and then we find conditions to characterize which one is better. Recall that $I' = q \left((1-p)I + p\mathbb{E}_I^{(f)} \right)$

4.4.1 With checkpoints during prediction window (Algorithm 1)

In order to compute the optimal value for T_P , let us find the portion of the waste that depends on T_P :

$$\text{WASTE}_{T_P} = \frac{rq}{\mu} \left(\frac{(1-p)I + p\mathbb{E}_I^{(f)}}{p} \frac{C}{T_P} + T_P \right)$$

As we can see, the optimal value for T_P is independent from q , but also from μ . The optimal value for T_P is thus:

$$T_P^{\text{extr}} = \sqrt{\frac{(1-p)I + p\mathbb{E}_I^{(f)}}{p}}C \quad (10)$$

However, for our algorithm to be correct, we want $\frac{I}{T_P} \in \mathbb{N}$ (the interval I is partitioned in k intervals of length T_P , for some integer k). We choose T_P^{opt} equal to either $\left\lceil \frac{I}{T_P^{\text{extr}}} \right\rceil$ or $\left\lceil \frac{I}{T_P^{\text{extr}}} \right\rceil + 1$, depending on the value that minimizes WASTE_{T_P} . Note that we also have the constraint $T_P^{\text{opt}} \geq C$, hence if both values are lower than C , then $T_P^{\text{opt}} = C$.

Now that we know that T_P^{opt} is independent from both q and T_{NP} , we can see the waste in Equation (6) as a function of two variables. One can see from Equation (6) that the waste is an affine function of q . This means that the minimum is always reached for either $q = 0$ or $q = 1$. We now consider the two functions $\text{WASTE}_{\text{withCkpt}\{q=0\}}$ and $\text{WASTE}_{\text{withCkpt}\{q=1\}}$ in order to minimize them with respect to T_{NP} . First we have:

$$\text{WASTE}_{\text{withCkpt}\{q=0\}} = \frac{C}{T_{NP}} + \frac{1}{\mu} \left(\frac{T_{NP}}{2} + D + R \right) \quad (11)$$

As expected, this is exactly the equation without prediction, the study of the optimal solution has been done in Section 3, it is minimized when $T_{NP}^{\text{opt}_0} = \min(\alpha\mu_e - I, \max(\sqrt{2C\mu}, C))$.

Next we have:

$$\begin{aligned} \text{WASTE}_{\text{withCkpt}\{q=1\}} &= \left(1 - \frac{r((1-p)I + p\mathbb{E}_I^{(f)})}{p\mu} \right) \left(\frac{C}{T_{NP}} + \frac{1-r}{\mu} \frac{T_{NP}}{2} \right) \\ &+ \frac{r}{\mu} \left(\frac{((1-p)I + p\mathbb{E}_I^{(f)})}{p} \frac{C}{T_P^{\text{opt}}} + T_P^{\text{opt}} \right) + \frac{r}{p\mu} C \\ &+ \left(\frac{r}{\mu} + \left(1 - \frac{r((1-p)I + p\mathbb{E}_I^{(f)})}{p\mu} \right) \frac{1-r}{\mu} \right) (D + R) \end{aligned} \quad (12)$$

This equation is minimized when

$$T_{NP}^{\text{opt}_1} = \sqrt{\frac{2\mu C}{(1-r)}}$$

One can remark that this value is equal to the result without intervals (Section 3). Actually, the only impact of the prediction interval I is the moment when we should take a pre-emptive action. Note that when $r = 0$ (this means that there is no prediction), we have $T_{NP}^{\text{opt}_1} = T_{NP}^{\text{opt}_0}$, and we retrieve Young's formula [4].

Finally, we know that the waste is defined for $C \leq T_{NP} \leq \alpha\mu_e - I$. Hence, if $T_{NP}^{\text{opt}_1} \notin [C, \alpha\mu_e - I]$, this solution is not satisfiable. However Equation (12) is convex, so the optimal solution is C if $T_{NP}^{\text{opt}_1} < C$, and $\alpha\mu_e - I$ if $T_{NP}^{\text{opt}_1} > \alpha\mu_e - I$. Hence, when $q = 1$, the optimal solution should be

$$\min \left(\alpha\mu_e - I, \max \left(\sqrt{\frac{2\mu C}{(1-r)}}, C \right) \right). \quad (13)$$

4.4.2 Waste for the algorithm that does not checkpoint during the proactive mode

One can see that Equation (8) and Equation (6) only differ: by the quantity $\frac{qr}{\mu} \left(\frac{(1-p)I + p\mathbb{E}_I^{(f)}}{p} \frac{C}{T_P^{\text{opt}}} + T_P^{\text{opt}} - \mathbb{E}_I^{(f)} \right)$. This value is linear in q and a constant with regards to T_{NP} . Hence the minimization is almost the same.

Once again we can see that the optimal value for q is either 0 or 1. We can consider the two functions $\text{WASTE}_{\text{noCkpt}\{q=0\}}$ and $\text{WASTE}_{\text{noCkpt}\{q=1\}}$. We remark that $\text{WASTE}_{\text{noCkpt}\{q=0\}} = \text{WASTE}_{\text{withCkpt}\{q=0\}}$,

and hence that the study has already been done. As for $\text{WASTE}_{\text{noCkpt}\{q=1\}}$, it is also minimized when

$$T_{\text{NP}}^{\text{opt}} = \sqrt{\frac{2\mu C}{(1-r)}}.$$

Finally, the last step of this study is identical to the previous minimization, and the optimal solution when $q = 1$ is defined by $T_{\text{NP}}^{\text{opt}_1} = \min \left(\alpha\mu_e - I, \max \left(\sqrt{\frac{2\mu C}{(1-r)}}, C \right) \right)$.

4.4.3 Identifying the most efficient algorithm

Finally in this section, we consider the waste for the two algorithms that take the prediction window into account (the one that does not checkpoint during the prediction window, and the one that checkpoints during the prediction window), and try to find conditions of dominance of one strategy over the other. Since the equation of the waste is identical when $q = 0$, let us consider the case when $q = 1$. We have seen that:

$$\text{WASTE}_{\text{withCkpt}\{q=1\}} - \text{WASTE}_{\text{noCkpt}\{q=1\}} = \frac{r \left((1-p)I + p\mathbb{E}_I^{(f)} \right)}{p\mu} \frac{C}{T_{\text{P}}^{\text{opt}}} + \frac{r}{\mu} \left(T_{\text{P}}^{\text{opt}} - \mathbb{E}_I^{(f)} \right) \quad (14)$$

We want to know when Equation (14) is nonnegative (meaning that it is beneficial not to take any checkpoints during proactive mode). We know that this value is minimized when $T_{\text{P}}^{\text{extr}} = \sqrt{\frac{(1-p)I + p\mathbb{E}_I^{(f)}}{p}} C$ (Equation (10)), then a sufficient condition would be to study the equation $\text{WASTE}_{\text{withCkpt}\{q=1\}} - \text{WASTE}_{\text{noCkpt}\{q=1\}} \geq 0$ with $T_{\text{P}}^{\text{extr}}$ instead of $T_{\text{P}}^{\text{opt}}$. That is:

$$\begin{aligned} \frac{r(1-p)I + p\mathbb{E}_I^{(f)}}{p\mu} \frac{C}{\sqrt{\frac{(1-p)I + p\mathbb{E}_I^{(f)}}{p}} C} + \frac{r}{\mu} \left(\sqrt{\frac{(1-p)I + p\mathbb{E}_I^{(f)}}{p}} C - \mathbb{E}_I^{(f)} \right) &\geq 0 \\ \Leftrightarrow 2\sqrt{\frac{(1-p)I + p\mathbb{E}_I^{(f)}}{p}} C &\geq \mathbb{E}_I^{(f)} \end{aligned} \quad (15)$$

Consequently, we can say that if Equation (15) is matched, then $\text{WASTE}_{\text{noCkpt}} \leq \text{WASTE}$, the algorithm where we do not checkpoint during the proactive mode has a better solution than Algorithm 1. For example, if we assume that faults strike uniformly during the prediction window $[t_0, t_0 + I]$, in other words, if $0 \leq x \leq I$, the probability that the fault occurs in the interval $[t_0, t_0 + x]$ is $\frac{x}{I}$, then $\mathbb{E}_I^{(f)} = \frac{I}{2}$, and our condition becomes

$$I \leq 16 \frac{1-p/2}{p} C.$$

We can now finish our study by saying that in order to find the optimal solution, one should compute both optimal solutions for $q = 0$ and $q = 1$, for both algorithms, and choose the one that minimizes the waste, as was done in Section 3, except when Equation (15) is valid, then we can focus on the computation of the waste of the algorithms that does not checkpoint during proactive mode.

5 Related work

Considerable research has been conducted on fault prediction using different models (system logs analysis, event-driven approach). In this section we give a brief overview of the results obtained by predictors. We focus on their results rather than on their methods of prediction.

The authors of [6] introduce the *lead time*, that is the time between the prediction and the actual fault. It is a time that should be sufficient to take proactive actions. They are also able to give the location of the fault. The fact that the location is given has an impact on the precision and the recall (see Table 5). The authors of [5] also consider also a lead time, and introduce a *prediction window* when the predicted fault should happen. This motivates the work on Section 4, even though they never give the size of their prediction window. Unfortunately, much of the work done on prediction does not

Paper	Lead Time	Precision	Recall	Prediction Window
[6]	300 s	40 %	70%	-
[6]	600 s	35 %	60%	-
[5]	2h	64.8 %	65.2%	yes (size unknown)
[5]	0 min	82.3 %	85.4 %	yes (size unknown)
[1]	32 s	93 %	43 %	-
[3]	NC	70 %	75 %	-

Table 1: Comparative study of different parameters returned by some predictors.

provide information that could be really useful for the design of good algorithms. These informations are those stated above, namely the lead time and the size of the prediction window, but other information that could be useful would be the distribution of the faults in the prediction window, the precision as a function of the recall (see our analysis), or even the precision and recall as functions of the prediction window (what happens with a bigger prediction window). Again, all these informations could be useful in the design of good algorithms.

While many study on fault prediction focus on the conception of the predictor, most of them consider that the proactive action should simply be a checkpoint of a migration right before the fault. However, in their paper [2], Li et al. consider the mathematical problem to determine when and how to migrate. In order to be able to use migration, they stated that at every time, 2% of the resources are available. This allowed them to conceive a Knapsack-based heuristic. Thanks to their algorithm, they were able to save 30% of the execution time compared to an heuristic that does not take the reliability into account, with a precision and recall of 70%, and with a maximum load of 0.7.

6 Conclusion

The comprehensive analytical results provided in this preliminary report enable to fully asses the impact of fault prediction on optimal checkpointing strategies. Future work will be devoted to instantiate these results with realistic application/platform scenarios, and to provide an experimental evaluation of the importance of fault prediction to reduce checkpoint overhead.

Acknowledgments. The authors are with Université de Lyon, France. Y. Robert is with the Institut Universitaire de France. This work was supported in part by the ANR *RESCUE* project.

References

- [1] B. Kramer A. Gainaru, F. Cappello. Taming of the shrew: Modeling the normal and faulty behavior of large-scale hpc systems. In *Proceedings of the 2012 IEEE International Parallel and Distributed Processing Symposium*, IPDPS’12, may 2012.
- [2] Song Fu and Cheng-Zhong Xu. Exploring event correlation for failure prediction in coalitions of clusters. In *Proceedings of the 2007 ACM/IEEE conference on Supercomputing*, SC ’07, pages 41:1–41:12, New York, NY, USA, 2007. ACM.
- [3] Errin W. Fulp, Glenn A. Fink, and Jereme N. Haack. Predicting computer system failures using support vector machines. In *Proceedings of the First USENIX conference on Analysis of system logs*, WASL’08, pages 5–5, Berkeley, CA, USA, 2008. USENIX Association.
- [4] John W. Young. A first order approximation to the optimum checkpoint interval. *Comm. of the ACM*, 17(9):530–531, 1974.
- [5] Li Yu, Ziming Zheng, Zhiling Lan, and S. Coghlan. Practical online failure prediction for blue gene/p: Period-based vs event-driven. In *Dependable Systems and Networks Workshops (DSN-W), 2011 IEEE/IFIP 41st International Conference on*, pages 259 –264, june 2011.

- [6] Ziming Zheng, Zhiling Lan, R. Gupta, S. Coghlan, and P. Beckman. A practical failure prediction with location and lead time for blue gene/p. In *Dependable Systems and Networks Workshops (DSN-W), 2010 International Conference on*, pages 15 –22, 28 2010-july 1 2010.